

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 04-19177

(43)Date of publication of application : 10.07.1992

(51)Int.Cl.

G06F 15/40

G06F 15/18

(21)Application number : 02-323540

(71)Applicant : NIPPON TELEGR & TELEPH CORP <NTT>

(22)Date of filing : 27.11.1990

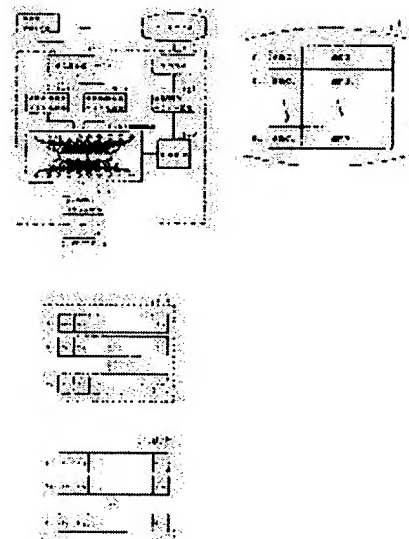
(72)Inventor : KUWABARA SATOSHI

(54) SIMILARITY CORRECTION PROCESSING METHOD

(57)Abstract:

PURPOSE: To find out a program with a high utility value by retrieving similar data reflecting an experienced retrieval tendency when a characteristic with a constitution is turned to be a retrieval key and the similarity retrieval of the similar data is performed based on a data base.

CONSTITUTION: Routing St_i is calculated between retrieval data C_t and the whole characteristic data C_i ($i=1$ to N) in the data base by a similarity generation part 3-1, one N -dimensional similarity vector $St=(st_1, st_2, \dots, st_N)$ is prepared to be stored in similarity vector storage part 3-3 for retrieval. Corresponding rank vector $Et=(et_1, et_2, \dots, et_N)$ is outputted to a rank vector storage part 3-8 for retrieval by inputting this similarity vector in a learned neural circuit network 3-6. In this rank vector, the rank to be displayed indicates the rank of a similar data pair to be retrieved actually, and when the learning in the former stage is sufficiently refined, this rank reflects the experienced retrieval tendency. Thus, the program with high utility value can be find out.



LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

⑫ 公開特許公報(A) 平4-191977

⑤ Int. Cl.⁵G 06 F 15/40
15/18

識別記号

5 1 0 J

庁内整理番号

7056-5L
8945-5L

⑬ 公開 平成4年(1992)7月10日

審査請求 未請求 請求項の数 1 (全9頁)

⑭ 発明の名称 類似度補正処理方法

⑯ 特 願 平2-323540

⑰ 出 願 平2(1990)11月27日

⑱ 発 明 者 桑 原 敏 東京都千代田区内幸町1丁目1番6号 日本電信電話株式会社内

⑲ 出 願 人 日本電信電話株式会社 東京都千代田区内幸町1丁目1番6号

⑳ 代 理 人 弁理士 森 田 寛

明 細 書

1. 発明の名称

類似度補正処理方法

2. 特許請求の範囲

特徴データと操作データとからなるN個のデータ対が格納されたデータベースから、検索データの特徴を検索キーに、それと類似な特徴を有するデータ対を検索する際に、

2つの特徴データ間の類似度を生成する類似度生成手段と、

前記類似度生成手段を用いて作成した1個の検索データとN個のデータベース内のデータ対とのN次元類似度ベクトルを記憶する検索用類似度ベクトル記憶手段と、

前記類似度生成手段を用いて作成したデータベース内のN個のデータ対相互間のN×N次元類似度ベクトルを記憶する学習用類似度ベクトル記憶手段と、

データベース内のデータ対相互間のN×N次元
序列ベクトルを生成する序列生成手段と、

前記N×N次元序列ベクトルを記憶する学習用
序列ベクトル記憶手段と、

1個のN次元検索用類似度ベクトル、あるいは
1個のN次元学習用類似度ベクトルを入力とし、
1個のN次元序列ベクトルを出力とする神経回路
網とを設け、

検索に先だって、予め前記学習用類似度ベクトル
のうち1個を神経回路網に入力し、教師信号を
上記学習用類似度ベクトルに対応した1個の学習
用序列ベクトルとし、それぞれ学習用に蓄積され
たN個のベクトルについて学習を繰り返すことに
より、神経回路網の最適化を行い、

検索データが与えられた時は、検索用類似度ベ
クトルを前記学習済みの神経回路網へ入力するこ
とにより、実際に有益なデータの序列を示す序列
ベクトルを得るようにした

ことを特徴とする類似度補正処理方法、

3. 発明の詳細な説明

(産業上の利用分野)

本発明は、ある構造を持った特徴を検索キーにして、それに類似なデータをデータベースから検索する(以下類似検索)というシステムにおいて、特徴間の構造的類似関係によってのみ求められる類似度に基づいて類似検索するのではなく、経験的な検索傾向による類似度の補正を行い、より実用的な類似なデータ検索を可能にするようにした類似度補正処理方法に関するものである。

例えば類推を用いた化学物質設計支援システムにおいて類似検索を応用した場合、化学物質の反応事例を反応後の化学構造を特徴データとし、反応前の化学構造及び反応条件を操作データとして、多数格納された化学物質反応データベースから、設計対象とする化学物質の化学構造を特徴とする検索キーにして、それと類似した化学物質の反応データをデータベースから検索する必要がある。その場合の設計の良否は、より実用的なデータを検索することである。しかし構造特徴のような

正し、この補正された類似度の序列から、より実用性の高いデータを検索することが考えられてきた。さらに本発明との係わりでは、この重み係数を、データベースから帰納的に決定する方法が提案されてきた(情報処理学会第40回(平成2年度前期)全国大会講演論文集(1)4D-6「重みを利用した類似化学反応の検索」による)。この従来技術における重みの帰納法的決定方法とは、データベースから単一のデータを引き抜き、そのデータとその他のデータとの間の類似度値の序列を求め、さらにそのデータとその他のデータとの間の実際の有用な序列を類推シミュレーションによって求め、両者の序列間の差分を検出し、最初に求められた類似度序列のうち過小に見積られた類似度を大きくするように、対応するデータに付けられた重みを単位置だけ増加することを、全てのデータについて繰り返し行う方法と考えられる。

(発明が解決しようとする課題)

上記の従来技術方式における重み決定方法は、

面的な類似性だけにより、類似なデータを検索することは必ずしも、設計に携わる利用者にとって実用的なデータとは限らない。本発明はこのような類推を用いた設計支援システム等に必要な類似検索に関して、検索側のデータ特徴とデータベース内の特徴データとの構造的な最大類似性のみを考慮して、類似データ対を検索するのではなく、検索されようとしているデータ対の利用頻度や重要性等の経験的な検索傾向をも考慮した類似検索を可能にし、より実用性の高い設計支援システム等を実現可能とすることを狙いとしている。

(従来技術)

初期の類似検索では、特徴の構造的類似性を部分種大マッチング等により、定量的な類似度で表現し、この類似度の序列でデータ検索を行う方法が採用されてきた。さらに、改良方法として、構造的な類似度の序列が実際に有用なデータの順序を保証しているわけではないことから、個々のデータに重み係数を付与し、これにより類似度を補

事例に基づいた一種の強化学習を行うのであるが、問題点としては

- ① 修正が関係するデータに付与された重みに、単位置でしか行われぬこと、
- ② 重み修正が増加方向にしかなされないことのために、必ずしも重み決定の収束保証がないこと

があげられ、これが類似検索の精度を低下させていると考えられる。

本発明は上記従来技術の問題点に関し、①については、関係するデータの重みだけでなく全データの重みに対して、可変量で修正するようにしたこと、また②については、重みの修正を増加方向だけでなく、減少方向にも行うようにしたことであり、これにより重み決定の収束性を向上するとともに、類似度補正の精度を向上することを目的としている。

(課題を解決するための手段)

本発明は類似検索における構造的類似度の序列

を実際に有用な類似度の序列にマッピングするための類似度補正に関するものである。これは、データベース中の全データ数 (N 個) を次元数とする N 次元の類似度ベクトルを同じ N 次元の序列ベクトルに変換することを意味する。この変換機能は、変換事例が類似度ベクトルと序列ベクトルとして獲得できる場合には、事例学習によって構築することができる。本発明では、事例学習の手段として、 $N \times N$ の神経回路網を利用することとし、学習に必要な事例については、データベース中のデータから獲得することとする。つまり、データベースからデータを1つ抜き出し、そのデータとデータベース中の N 個のデータとのそれぞれの類似度を計算し、それを N 次元ベクトル化したものを入力信号、またそのデータと N 個のデータとのそれぞれの序列を類推によって推定し、それを N 次元ベクトル化したものを教師信号とし、これらを組にして1つの事例データとする。従って、 N 個のデータについて N 組の事例データが得られ、学習段階ではこれらを繰り返し用いることとする。

(実施例)

本発明に係わる機能構成について第1図に実施例を示す。

第1図においては、検索用特徴データ1、データベース2、序列変換部3、及び序列データ4から構成される。さらに、第2図は序列変換部の内部構成を示しており、1、2、3、4は第1図に対応し、類似度生成部3-1、序列生成部3-2、検索用類似度ベクトル記憶部3-3、学習用類似度ベクトル記憶部3-4、学習用序列ベクトル記憶部3-5、及び神経回路網3-6、学習機構3-7、及び検索用序列ベクトル記憶部3-8から構成される。

なお、検索用特徴データはデータベース内データ対の特徴データと同一形式であり、検索用類似度ベクトルは学習用類似度ベクトルの1エントリと同一形式であり、検索用序列ベクトルは学習用序列ベクトルの1エントリと同一形式である。

序列データ4は検索用序列ベクトルのサブセットであり、実際に必要なデータ対数だけを提示す

(作用)

データベースのデータを用いて学習するために同じデータを使った神経回路網による想起が可能であるばかりでなく、データベースに無いデータについても近似的な連想機能が利用できる。つまり、新規のデータに対して、全データとの計算された N 次元類似度ベクトルを入力とした時に、出力には N 次元の実際に有効な序列ベクトルが得られる。

前述のことは神経回路網を使用することにより、従来技術において個々のデータに付与していた重みは神経回路網内部のユニット間の重みに置き換えられ、重みによる類似度補正を神経回路網内のユニット間の重み計算とユニット内の伝達関数の働きによって達成したことであり、さらに収束性の保証されたバックプロパゲーション等の神経回路網を利用した事例学習により、より正確な類似度補正を可能にしたことである。

るために設けられる。但し、この機構については本発明の請求範囲外のことであり詳細は省く。

本発明の動作について説明する。動作は大きく学習フェーズと検索フェーズとから構成され、学習フェーズが検索フェーズに先行する。即ち初期の神経回路網3-6はランダムに初期化されており、実際の検索に先だっては適切な動作が可能ないように事例からの神経回路網3-6の重み設定が必要である。これを学習といい、検索段階に先だてて行われる必要がある。

以降、学習段階、検索段階の順に実施例を説明する。第8図は本発明の実施例概略機能フローチャートを表す。

(学習段階)

データベース2は第3図図示のように特徴データ C_i と操作データ O_i との対からなる N 個のデータ対 D_i ($i=1 \sim N$) が格納されているものとする。学習段階では、先ず学習用類似度ベクトル記憶部3-4と学習用序列ベクトル記憶部3-5とを準備する。学習用類似度ベクトル記憶部3

-4 はデータベース 2 内の N 個のデータ対 D_i ($i = 1 \sim N$) 相互間の類似度 s_{ij} ($i = 1 \sim N, j = 1 \sim N$) を類似度生成部 3-1 により計算し、第 4 図に示すように $N \times N$ ベクトルとして記憶したものである。また学習用序列ベクトル記憶部 3-5 は学習用類似度ベクトル記憶部 3-4 に対応して、データベース内の N 個のデータ対 D_i ($i = 1 \sim N$) 相互間の序列 e_{ij} ($i = 1 \sim N, j = 1 \sim N$) を序列生成部 3-2 により計算し、第 5 図に示すように $N \times N$ ベクトルとして記憶したものである。類似度生成部 3-1 は 2 つの特徴データ C_p, C_q 間の類似度 s_{pq} を生成するものであり、化学構造図のように特徴データがグラフ構造式で記述されておれば、グラフ間の最大一致度を類似度とすることができる。さもないとすれば 2 つの特徴データからそれぞれ同一のカテゴリで特徴表現しベクトル化した 2 つの特徴ベクトル間での一致数を類似度とすることもできる。いずれにしても類似度は対象とする問題領域毎に適切に作成されたものを利用するものとする。

類推については対象とする問題に依存するものである。以下化学反応を例とした場合の類似度生成、及び序列生成について説明する。

第 6 図は類似度生成部 3-1 の働きを示したものである。2 つの特徴データ C_p, C_q が示されている。特徴データ C_p, C_q は化学構造を特徴インデックスとして捉えた表現が用いられており、類似度 s_{pq} は特徴データ内のインデックス間での一致度の総得点としている。あるいはここで各インデックス間の一致度に適当に重みを付加した加重加算による総得点とすることもできる。即ち、特徴を $C_p = (c_{p1}, c_{p2}, \dots, c_{pk})$ と $C_q = (c_{q1}, c_{q2}, \dots, c_{qk})$ とした場合、類似度 s_{pq} は以下のように計算される。

$$s_{pq} = \sum_k a_k \{c_{pk}, c_{qk}\}$$

但し、 $\{c_{pk}, c_{qk}\}$: 特徴インデックス k 間の一致度

a_k : 特徴インデックス k の重み係数

例えば、第 6 図では 2 種の芳香族単環化学物質

序列生成部 3-2 は 2 つのデータ対 D_p, D_q 間の序列度を生成するものであり、化学反応のような構造変化に対してはデータベース自身から類推シミュレーションを利用して生成することができる。即ちデータ対 D_p について特徴データ C_p 、特徴データ C_q 、及び操作データ O_q を利用して操作データ O_p を類推する類推手段を有することが前提である。1 つのデータ対 D_i 内の特徴データ C_i についてデータベース内の全てのデータ対 D_j との間で操作データを類推すると、 N 個の操作データ O_{ij} ($j = 1 \sim N$) が得られる。この N 個の操作データ O_{ij} ($j = 1 \sim N$) をデータ対 D_i 内の既知である操作データ O_i と比較し、類似度の大きいものから順番に 1 から最大 N までの序列を付ける。このデータ対 D_i のデータ対 D_j に対する序列を e_{ij} ($j = 1 \sim N$) とし、全ての i ($i = 1 \sim N$) について序列 e_{ij} ($i = 1 \sim N, j = 1 \sim N$) を求めれば $N \times N$ のベクトルが生成される。但し、類推が不能な場合は序列の値は '0' とする。なお、この序列化に用いる

間の類似度を化学構造特徴インデックスを適当に設定した場合に、類似度が 2 2 と計算されることを示している。このような計算手法を用いて、あるデータ対 D_i 内の特徴データ C_i について他のデータ対 D_j 内の特徴データ C_j との類似度 s_{ij} を求め、 N 次元の類似度ベクトル $S_i = (s_{i1}, s_{i2}, \dots, s_{iN})$ を作成する。さらに全ての特徴データ C_i について同様に N 個の N 次元類似度ベクトル S_i ($i = 1 \sim N$) が生成でき、これを学習のために使用する学習用類似度ベクトル記憶部 3-4 に格納する。

次に序列生成部 3-2 の働きについて説明する。

化学反応は化学構造の変化であり、その変化の前後がそれぞれデータ対内の特徴データ C 、及び操作データ O に格納されているものとする。同種の化学変化を持つ複数のデータ対相互では化学構造変化が類似しており、一方の変化前あるいは変化後の化学構造が無くても類推により生成可能である場合が多い。この類推の手法があるものとする。2 つのデータ対 D_p, D_q について一方

の特徴データ C_p のみと他方のデータ対 (C_q, O_q) とから操作データ O_{pq} を類推し、これが既知の操作データ O_p と類似している場合、この2つのデータ対 D_p と D_q は類似していると考えられる。この類似性は構造的な類似度 a_{pq} として容易に生成可能であり、先の類似度生成機構を使用しても良い。つまり、 a_{pq} は以下のように定義できる。

$$a_{pq} = \{O_p, O_{pq}\} = \{O_p, (C_p, C_q, O_q)\}$$

但し、 $\{O_p, O_{pq}\}$: 操作データ間の類似度

$$O_{pq} = (C_p, C_q, O_q)$$

(C_p, C_q, O_q) : O_{pq} の類推

1つのデータ対 D_i について他のデータ対 D_j ($j = 1 \sim N$) との類似度 a_{ij} を求め、1つの類似度ベクトル $A_i = (a_{i1}, a_{i2}, \dots, a_{iN})$ が求められるが、本発明ではこれを1からNの序列値で序列化したものを序列ベクトル $E_i = (e_{i1}, e_{i2}, \dots, e_{iN})$ とする。つまり各データ対 D_i について、1個のN次元序列ベクトル E_i が生成でき、全データ対 D_i については $N \times N$ 次元の序列

ベクトル E_i ($i = 1 \sim N$) が生成でき、これを学習のために使用する学習用序列ベクトル記憶部3-5に格納する。なお、類推時に類推が不能な場合がある。例えば、2つのデータ対 D_i, D_j 間に同一化学反応を持つことができない場合、即ち反応中心を共有できない場合は類推するまでもなく、類似度 a_{ij} は '0' であり、序列値 e_{ij} も '0' とみなされる。

第7図は4つのデータ対 D_p, D_q, D_r 、及び D_v について D_p と D_q, D_r 、及び D_v との序列生成例を示したものである。なお、ここでは D_v のように反応中心を共有できない場合や、 D_r のように類推によって生成した O_{pr} と O_p が一致しない場合には類似度 a_{pr} 、 a_{pv} は '0' とし、 D_q の場合のように類推によって生成した O_{pq} と O_p が一致する場合には a_{pq} は '1' とした。この場合には類似度 a_{pj} ($j = q, r, v$) はそのまま序列 e_{pj} ($j = q, r, v$) として使用される。

次に、神経回路網の学習について説明する。神経回路網3-6を初期化後、学習用類似度ベクトル

記憶部3-4から $i = 1$ の前期類似度ベクトル A_i が神経回路網3-6の入力に加えられ、学習機構3-7へは教師信号として学習用序列ベクトル記憶部3-5から $i = 1$ の序列ベクトル E_i が与えられる。学習機構では神経回路網の出力である序列ベクトル E_i と学習機構に与えられた予め計算済みの序列ベクトル E_i とから誤差量を計算する。誤差量はバックプロパゲーション学習アルゴリズムで最も一般的な二乗誤差等が用いられる。学習機構3-7ではこの誤差の変化を基にバックプロパゲーション学習を行い、神経回路網の重み等を修正する。引き続き、学習サイクルを $i = 2$ から $i = N$ まで行い、さらにこれを誤差が一定の値まで降下するまで同一学習用類似度ベクトル記憶部3-4、及び学習用序列ベクトル記憶部3-5内のそれぞれ類似度ベクトル、及び序列ベクトルを複数回用いて学習サイクルを繰り返す。以上により、出力の序列ベクトルと期待される序列ベクトルとの誤差がある一定量以下になった場合に神経回路網3-6の学習は終了したものとみなし、

学習フェーズを終了する。

(検索段階)

検索フェーズではデータベース内の特徴データと同じ形式の検索データ C_t が与えられたとした場合には、類似度生成部3-1の働きにより、検索データ C_t とデータベース内の全特徴データ C_i ($i = 1 \sim N$) 間で類似度 S_{ti} が計算され1つのN次元類似度ベクトル $S_t = (s_{t1}, s_{t2}, \dots, s_{tN})$ が作成され、検索用類似度ベクトル記憶部3-3に記憶される。この類似度ベクトルを学習済みの神経回路網3-6に入力すると、対応する序列ベクトル $E_t = (e_{t1}, e_{t2}, \dots, e_{tN})$ を検索用序列ベクトル記憶部3-8に出力する。この序列ベクトルにおいて、表示される序列が実際に検索されるべき類似なデータ対の序列を示しており、前段の学習が充分精練されているならば、この序列は経験的な検索傾向を反映するはずであり、実際には、そのうちの上位いくつかを序列データ4として出力する。この序列データが示すデータベース内のデータ対を利用することによって、

本発明の目的とするより有効な設計が可能となる。

(発明の効果)

以上の説明から明らかになるように、本発明によれば、ある構造を持った特徴を検索キーにして、それに類似なデータをデータベースから類似検索する場合、特徴間の構造的類似関係によってのみ類似なデータを検索するのではなく、経験的な検索傾向をも反映した類似なデータ検索を行うことが可能になる。上記実施例では化学物質設計支援の場合について示したが、他に例えば、ソフトウェア設計支援システム等において類推手法によるプログラムモジュールの再利用を行う場合、単にプログラムの特徴が似ているということだけで類似なプログラムを検索するのではなく、プログラムの利用度、重要性、経験的な有効性を反映した類似度に補正でき、その結果、より利用価値の高いプログラムを見つけ出すことができる等、このような類似検索を必要とするシステムは多く考えられ、知識ベースから類似な知識を検索して類

推を行う他の一般的な知識ベースシステムについても同様な効果が期待できるという著しい工業的効果がある。

4. 図面の簡単な説明

第1図は類似度補正システムの機能構成、

第2図は序列変換部の機能構成、

第3図はデータベース内のデータ対構成、

第4図は学習用類似度ベクトル記憶部構成、

第5図は学習用序列ベクトル記憶部構成、

第6図は類似度生成部の働き、

第7図は序列生成部の働き、

第8図は類似度補正の概略機能フローである。

図中において、

1…検索用特徴データ

2…データベース

3…序列変換部

3-1…類似度生成部

3-2…序列生成部

3-3…検索用類似度ベクトル記憶部

3-4…学習用類似度ベクトル記憶部

3-5…学習用序列ベクトル記憶部

3-6…神経回路網

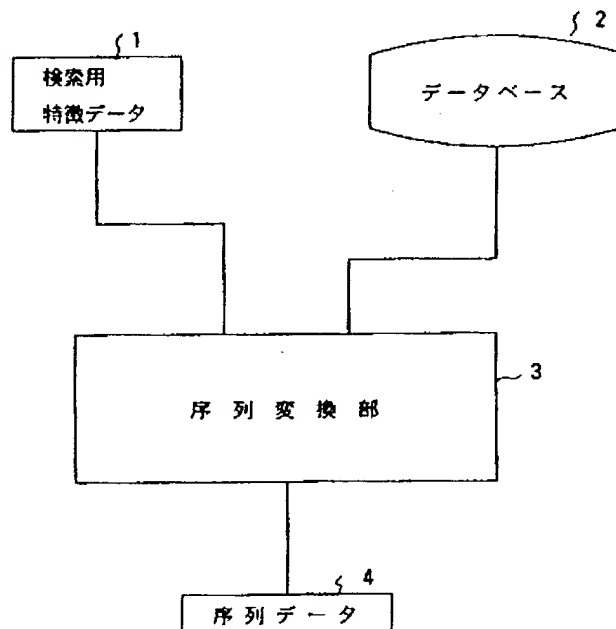
3-7…学習機構、

3-8…検索用序列ベクトル記憶部

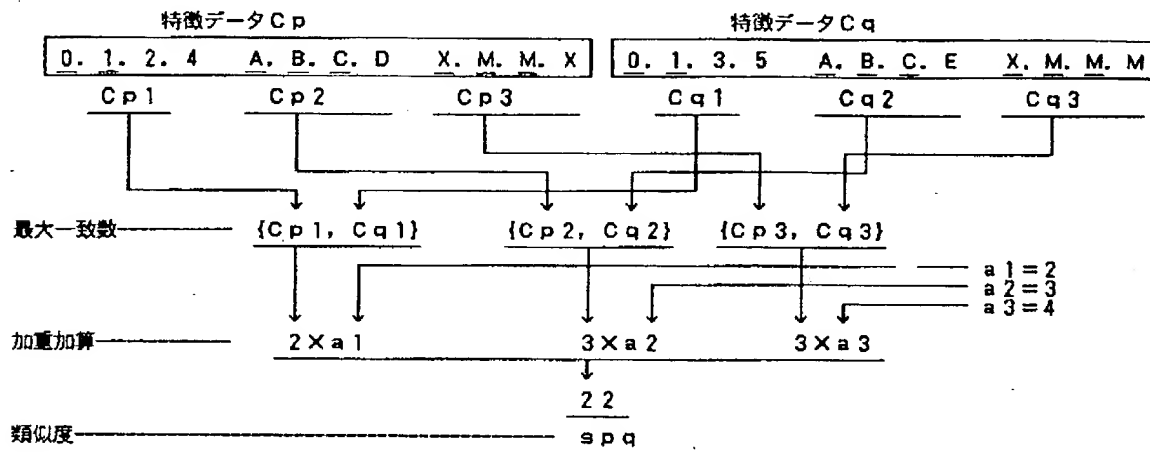
4…序列データ

特許出願人 日本電信電話株式会社

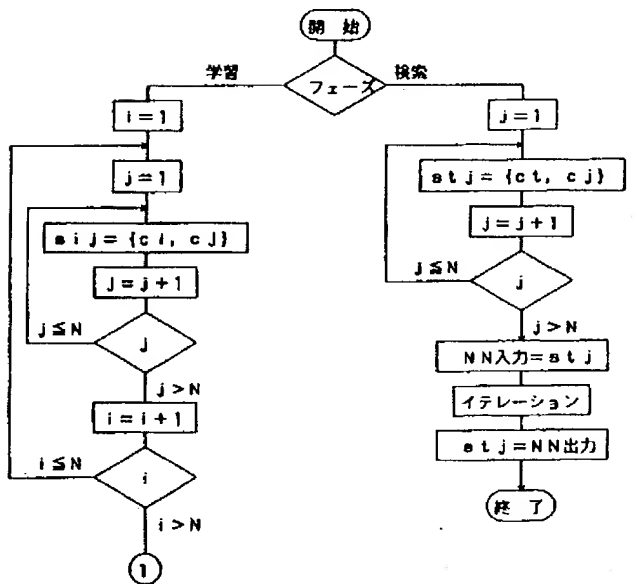
代理人 弁理士 森田 寛



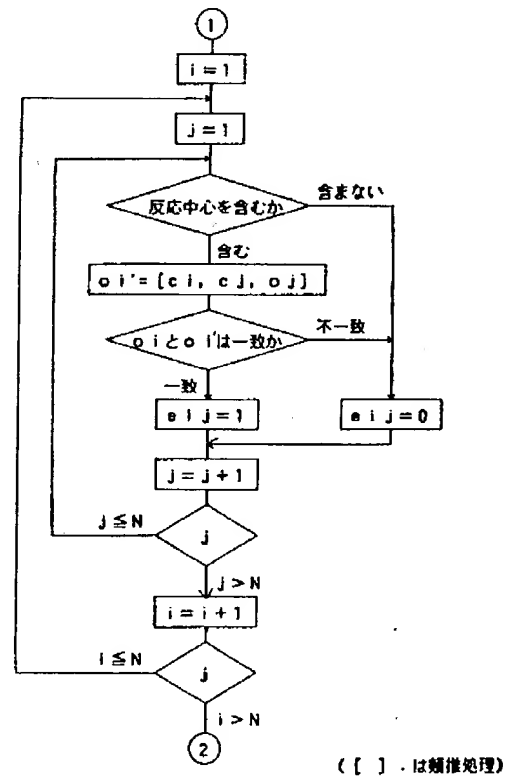
第1図



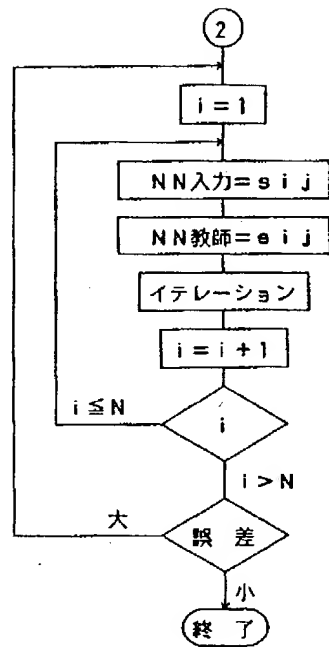
第 6 図



第 8 図 (I)



第 8 図 (II)



第 8 図 (Ⅲ)